# A FIRST ORDER ITERATION PROCESS FOR

## SIMULTANEOUS EQUATIONS*

H. A. Luther and W. F. Stewart[†]

## Introduction

The judgment of an iterative solution of simultaneous equations

properly entails more than the rate of convergence.  An item of con-

cern is whether or not we can arrive at a solution point whose

approximate location is known.

Let the "region of influence" of a given solution point be the

set of all starting values which (for a given iterative technique)

have that solution point as the limiting value of the iteration.

Ideally, perhaps, we might desire that every starting point yield

some solution and that every region of influence be simply connected.

It is known that, for example, the Newton-Raphson technique is

a second order process with favorable rate of convergence when well-

started.  It is also known that the regions of influence are not

normally simply-connected.

The present procedure was devised as part of a program for

studying methods of altering the regions of influence rather than

rates of convergence.

The equations considered can be linear or nonlinear.  By using

---

[†]Texas A&M University, College Station, Texas.

first order partial derivatives, there is created a technique which
has the Newton-Raphson quality of converging to every solution point
whose Jacobian is nonvanishing. The method is computationally
simpler than the Newton-Raphson method; it is, however, first order
rather than second order. An arbitrary function can be used para-
metrically.

The study is far from complete. The present approach seems
faulty in that in practice the convergence rate is too slow. Because
of this a method of accelerating convergence is included.

## Description of the Process

Let

(1) $\qquad f_i(x_1, x_2, \cdots x_n), \qquad 1 \le i \le n,$

be real functions defined for an open region $R$. Let
$X = [x_1 \; x_2 \; \cdots \; x_n]^t$ denote a column matrix, and write
$f_i(x_1, x_2, \cdots x_n) = f_i(X)$. Let $Y = [y_1 \; y_2 \; \cdots \; y_n]^t$ be a point of $R$
such that the following conditions are satisfied:

    a. For $1 \le i \le n$, $f_i(Y) = 0$ ;

    b. In a neighborhood of $Y$ the partial derivatives $\dfrac{\partial f_i}{\partial x_j}$ are

       continuous $(1 \le i, j \le n)$ ;

    c. For $f_{ij}(X) = \dfrac{\partial}{\partial x_j} f_i(X)$, the $n \times n$ matrix $J(X) = [f_{ij}(X)]$ is

       nonsingular at $Y$.

Let the eigenvalues of $\{1/Sp(J^t(Y) \; J(Y))\}J^t(Y) \; J(Y)$ be

$\lambda_1, \lambda_2, \cdots \lambda_n$. They are positive, and $\displaystyle\sum_{i=1}^{n} \lambda_i = 1$. Let them be so

ordered that $\lambda_i \geq \lambda_{i+1}$. Let d(X) be a function over R arbitrary except that it is continuous in a neighborhood of Y and except that

(2) $$0 < d(Y) < 2/\lambda_1 \ .$$

For $1 \leq i \leq n$ define functions $h_i(X)$ by

(3) $$h_i(x) = \{d(X) \sum_{\ell=1}^{n} f_\ell(X) \ f_{\ell i}(X)\}/Sp\big(J^t(X) \ J(X)\big) \ .$$

Let $H(X) = [h_1(X) \ h_2(X) \ \cdots \ h_n(X)]^t$. Let $X^{(k)} = [x_1^{(k)} \ x_2^{(k)} \ \cdots \ x_n^{(k)}]^t$.

The basic theorem can now be phrased as follows:

If $X^{(1)}$ is "near enough" Y and

(4) $$X^{(k+1)} = X^{(k)} - H\big(X^{(k)}\big)$$

then $\lim_{k \to \infty} X^{(k)} = Y$.

Finding suitable functions d(X) is not elaborate. Since $\lambda_1 < 1$ unless n = 1, we may for d(X) use any positive constant not greater than two. In the final stages of iteration $d(X)/Sp\big(J^t(X) \ J(X)\big)$ can be replaced by a constant. Indeed, for some systems it may be practical to choose d(X) as $2 \ Sp\big(J^t(X) \ J(X)\big)/Max_R\big(Sp(J^t(X) \ J(X))\big)$ from the start, thus simplifying computation.

A cautionary note seems in order. There can exist points Y, resulting from the limit process (4), such that J(Y) is singular. In that case it may not be that $f_i(Y) = 0$ for all i. Consider the example $f_1(x_1, x_2) = x_1 + x_2 - 1$, $f_2(x_1, x_2) = x_1 + x_2 - 2$. Then

$h_1 = h_2$ and, for $d(X) = 1$, is $x_1/2 + x_2/2 - 3/4$. For the starting values $x_1^{(1)} = \alpha$ and $x_2^{(1)} = \beta$, we find $x_1^{(k)} = \alpha/2 - \beta/2 + 3/4$ and $x_2^{(k)} = -\alpha/2 + \beta/2 + 3/4$ for $k \geq 2$.

## Linear Systems

We wish first to make an observation applicable to the general case. By Taylor's expansion

$$f_\ell(X) = \sum_{j=1}^{n} f_{\ell j}(\xi_\ell)(x_j - y_j)$$

where $\xi_\ell = [(y_1 + \overline{\theta x_1 - y_1}) \cdots (y_n + \overline{\theta x_n - y_n})]^t$. Now let

$$(5a) \qquad g_{ij}(X) = \{d(X)/Sp(J^t(X)\, J(X))\} \sum_{\ell=1}^{n} f_{\ell i}(X)\, f_{\ell j}(\xi_\ell)$$

and let

$$(5b) \qquad G(X) = [g_{ij}(X)], \qquad 1 \leq i,\, j \leq n \,.$$

Then clearly $\big(\text{see } (4)\big)$

$$(6) \qquad X - Y - H(X) = \big(I - G(X)\big)(X - Y) \,.$$

The only use we want of the above in the present section is to show that the eigenvalues $\mu_i$ of $I - G(Y)$ are given by $\mu_i = 1 - d(Y)\lambda_i$, while if (2) holds

$$(7) \qquad 1 > \mu_n, \quad \mu_{i+1} \geq \mu_i, \quad \mu_1 > -1 \,.$$

To see this, it is only necessary to observe that

(8) $$G(Y) = \{d(Y)/Sp[J^t(Y)\ J(Y)]\}J^t(Y)\ J(Y)$$

and that $d(Y) > 0$ together with $\lambda_i \geq \lambda_{i+1} > 0$ guarantee $\mu_{i+1} \geq \mu_i$

while $d(Y) < 2/\lambda_1$ means $\mu_1 > -1$ and $1 > \mu_n$.

Now for the linear case let

(9) $$f_i(X) = \sum_{j=1}^{n} a_{ij}\ x_j - b_i$$

and, for convenience only, consider $d(X) \equiv \delta$, $\delta$ being chosen so that it falls between zero and $2/\lambda_1$. Solution of (9) is of course solution of

(10) $$AY = B$$

where $A = [a_{ij}]$ and $B = [b_1\ b_2\ \circ\circ\circ\ b_n]^t$. Since $\frac{\partial}{\partial x_i} f_\ell(X) = a_{\ell i}$, it is a straightforward matter to see that

(11) $$X - H(X) = \left(I - \{\delta/Sp(A^tA)\}A^tA\right)X + \{d/Sp(A^tA)\}A^tB.$$

Thus convergence of process (4) is guaranteed for the linear case (see [1], pp. 161-170) since the eigenvalues of the coefficient of X in the right member of (11) are less in magnitude than 1.

It is true that if, say, $\delta = 1$, and $\delta/Sp(A^tA)$ is used as a constant multiplier, then even relatively large inaccuracies in the computation of $Sp(A^tA)$ have no influence on the theoretical accuracy of the iteration. However, it is also known that $A^tA$ is more ill-conditioned than $A$. This should be borne in mind in using the

process.

The Case n = 1. For n = 1, the process reduces to the solution
of f(y) = 0 by the technique

$$(12) \qquad x^{(k+1)} = x^{(k)} - d\left(x^{(k)}\right)f\left(x^{(k)}\right)/f'\left(x^{(k)}\right)$$

where d(y) = δ is a positive number less than two. For the choices
d(x) = 1, d(x) = $[f'(x)]^2/\{[f'(x)]^2 - f(x)f''(x)\}$ and related ideas
see [2], p. 24 et seq.

To establish convergence for a d(x) chosen in accordance with
the present discussion, let $R_1$ be an interval in R such that,
$R_1$, d(x) f'(y + $\overline{\theta x - y}$)/f'(x) = δ + ε(x) where θ is such that
f'(y + $\overline{\theta x - y}$)(x - y) is f(x) - f(y) and where $|\varepsilon(x)| < \mu - |1 - \delta|$
while 0 < μ < 1. Let $x^{(1)}$ be in $R_1$. Then $|x^{(k+1)} - y|$ is
$|x^{(k)} - y||1 - \delta - \varepsilon(x^{(k)})|$ which in turn does not exceed $\mu|x^{(k)} - y|$.
Thus $|x^{(k+1)} - y| \le \mu^k|x^{(1)} - y|$ and convergence occurs.

Now suppose d(y) = 1. Then convergence occurs as before, but
following a known technique (see [3], p. 448), can be shown to be
quadratic. Thus let d'(x) and f''(x) be continuous in $R_1$. Then,
expanding f(y) in powers of y - x, we are led to $x^{(k+1)} - y =$
$x^{(k)} - y - d\left(x^{(k)}\right)\{f'\left(x^{(k)}\right)(x^{(k)} - y) - f''\left(\eta^{(k)}\right)(x^{(k)} - y)^2/2\}/f'\left(x^{(k)}\right)$.
Then if $1 - d\left(x^{(k)}\right) = d(y) - d\left(x^{(k)}\right) = d'\left(\xi^{(k)}\right)\left(y - x^{(k)}\right)$, $x^{(k+1)} - y =$
$\left(x^{(k)} - y\right)^2\{d\left(x^{(k)}\right)f''\left(\eta^{(k)}\right) - 2d'\left(\xi^{(k)}\right)f'\left(x^{(k)}\right)\}/\{2f'\left(x^{(k)}\right)\}$.

The Case n = 2. In this instance the characteristic equation
for I - G(Y) is

knowledge of the functions is at hand.

Proof of Convergence.

Consider now $\bigl($see (4)$\bigr)$

$$X^{(k+1)} - Y = X^{(k)} - Y - H\bigl(X^{(k)}\bigr) .$$

Using (6) this becomes

$$X^{(k+1)} - Y = \bigl(I - G(X^{(k)})\bigr)\bigl(X^{(k)} - Y\bigr) .$$

Let $R_1$ be a spherical region in R with Y as center, and small enough that the functions $f_{ij}(X)$ and $d(X)$ are continuous in $R_1$. As in establishing (7), let the eigenvalues of $I - G(Y)$ be $\mu_i$, and let $\varepsilon$ be such that $1 - \varepsilon = \max_i(\mu_i^2)$. The numbers $\mu_i^2$ are the eigenvalues of $\bigl(1 - G(Y)\bigr)^t\bigl(I - G(Y)\bigr)$, and $1 - \varepsilon \geq \mu_i^2$. Let the eigenvalues of $\bigl(I - G(X)\bigr)^t\bigl(I - G(X)\bigr)$ be $\nu_i(X)$.

Let the characteristic polynomial of the $\mu_i^2$ be $\sum\limits_{i=0}^{n} a_i \mu^{2i}$ and that of the $\nu_i(X)$ be $\sum\limits_{i=0}^{n} b_i \nu^i$. Because $|b_i - a_i|$ can be made as small as desired by using a proper spherical subregion of $R_1$, there is in $R_1$ a spherical subregion $R_2$ having Y as its center and such that, for X in $R_2$, $|\nu_i - \mu_i^2| < \varepsilon/2$. In particular, $\max(\nu_i) < 1 - \varepsilon/2$. By a well-known theorem (see [1], p. 65)

$$\bigl(X^{(k)} - Y\bigr)^t\bigl(I - G(X)^{(k)}\bigr)^t\bigl(I - G(X^{(k)})\bigr)\bigl(X^{(k)} - Y\bigr) \leq$$

$$\max_i (\nu_i)\bigl(X^{(k)} - Y\bigr)^t\bigl(X^{(k)} - Y\bigr) \leq (1 - \varepsilon/2)\bigl(X^{(k)} - Y\bigr)^t\bigl(X^{(k)} - Y\bigr).$$

Then $\bigl(X^{(k+1)} - Y\bigr)^t\bigl(X^{(k+1)} - Y\bigr) \leq (1 - \varepsilon/2)\bigl(X^{(k)} - Y\bigr)^t\bigl(X^{(k)} - Y\bigr),$

and by recursion

$$||X^{(k+1)} - Y|| \leq (1 - \epsilon/2)^{k/2} ||X^{(1)} - Y||$$

Thus convergence occurs. Here $||X^{(k+1)} - Y||$ denotes the distance

norm.

Finally it is shown that under stringent conditions the convergence is quadratic. Let $J(Y)$ be a scalar multiple of an orthogonal matrix. Assume in addition that the first order derivatives of $d(X)$ and the second order derivatives of $f_\ell(X)$, $1 \leq \ell \leq n$, are continuous in a neighborhood of $Y$. Require also that $\delta = n$. Since in this instance all the eigenvalues $\lambda_i$ are equal, we have

$d(Y) = \delta = n < 2/\lambda_1 = 2n$. Thus convergence occurs and

$X^{(k)} - H(X^{(k)}) \big( \text{see } (4) \big)$ has meaning.

Write $d_i(X)$ for $\dfrac{\partial}{\partial x_i} d(X)$ and $f_{\ell ij}(X)$ for $\dfrac{\partial^2}{\partial x_i \partial x_j} f_\ell(X)$. Then

$$d(X) = \delta + \sum_{r=1}^{n} d_r(\xi)(x_r - y_r)$$

and, since $f_\ell(Y) = 0$

$$0 = f_\ell(X) + \sum_{r=1}^{n} f_{\ell r}(X)(y_r - x_r) + 1/2 \sum_{r,s=1}^{n} f_{\ell rs}(\eta)(x_r - y_r)(x_s - y_s).$$

First replace the functions $f_\ell(X)$ in $X - Y - H(X)$ by their equivalents above. The result, after dropping the quadratic terms, is

$X - Y - \{d(X)/Sp\big(J^t(X) \ J(X)\big)\}J^t(X) \ J(X) \ (X-Y)$. Now replace $d(X)$ by

$n + \sum_{r=1}^{n} d_r(\xi)(x_r - y_r)$. The result, after dropping the quadratic

terms, is

$$X - Y - \{n/\text{Sp}(J^t(X)\ J(X))\}J^t(X)\ J(X)\ (X - Y)\ .$$

Next in $J^t(X)\ J(X)$ replace, in the off diagonal terms, $f_{rs}(X)$ by $f_{rs}(Y) + \sum_{t=1}^{n} f_{rst}(\eta_{rs})(x_t - y_t)$. Because $J(Y)$ is a scalar multiple of an orthogonal matrix, there result only quadratic and cubic terms in $(x_i - y_i)$ for these components. There remains to consider

$$(15) \qquad X - Y - \{1/\text{Sp}(J^t(X)\ J(X))\}C(X)(X-Y)$$

where $C(X)$ is a diagonal matrix whose diagonal entry $c_i(X)$ is

$$n \sum_{\ell=1}^{n} f_{\ell i}^{2}(X)\ .$$

Now rephrase (15) as $\{1/\text{Sp}(J^t(X)\ J(X))\}(E(X) - C(X))(X - Y)$ where $E(X)$ is a diagonal matrix whose diagonal entry $e_i(X)$ is $\text{SP}(J^t(X)\ J(X))$. In $E(X) - C(X)$ replace each function $f_{ij}(X)$ by the Taylor's expansions used in arriving at (15). The result is a diagonal matrix which is made up of linear and quadratic terms in the $x_i - y_i$. This is so because $J(Y)$ is a scalar multiple of an orthogonal matrix.

In final consequence, $X - Y - H(X)$ is seen to contain only terms which are cubic or quadratic in the $x_i - y_i$.

## Improvement in Asymptotic Rate of Convergence

It was found in practice that process (4) converged slowly. Thus it was modified in an attempt to improve the asymptotic rate of

convergence.

    We start by defining

$$\overline{G}(X) \equiv \frac{d(X)}{Sp(J^t(X) \ J(X))} \ J^t(X) \ J(X)$$

where $J(X)$ and $Sp$ again denote the Jacobian matrix and spur respectively. The new process may now be stated as

$$(16) \qquad X^{(k+1)} = X^{(k)} - 2H(X^{(k)}) + \overline{G}(X^{(k)}) \ H(X^{(k)}) \ .$$

Then clearly

$$(17) \quad X - Y - 2H(X) + \overline{G}(X) \ H(X) = \left( I - 2G(X) + \overline{G}(X) \ G(X) \right) (X - Y)$$

where $G(X)$ is defined by (5).

    For the linear case, $\overline{G}(X) = G(X)$ and (17) becomes

$$(18) \qquad X - Y - 2H(X) + \overline{G}(X) \ H(X) = \left( I - G(X) \right)^2 (X - Y) \ .$$

As before let $d(X) \equiv \delta$ where $\delta$ is constant and consider a solution of (10). We have from (18)

$$(19) \qquad X - Y - 2H(X) + G(X) \ H(X) = (I - \frac{\delta}{Sp(A^t A)} A^t A)^2 \ X$$

$$+ \ (2I - \frac{\delta}{Sp(A^t A)} A^t A) \ \frac{\delta}{Sp(A^t A)} \ A^t B$$

Define M to be $(I - \frac{\delta}{Sp(A^t A)} A^t A)$ which is the coefficient of X in the right member of (11). Then $M^2$ is the coefficient of X in the right member of (19). If the eigenvalues of M are $\mu_i$, the eigenvalues of

$M^2$ are $\mu_i^2$. Since $|\mu_i| < 1$ (see (7)), then $\mu_i^2 < 1$ and we are guaranteed that process (16) converges for the linear case.

The asymptotic rate of convergence of M (see [4], p. 67) is defined as

(20)                           $R\infty\,(M) = -\ln\left(\rho(M)\right)$

where $\rho(M)$ is the magnitude of the largest eigenvalue of M. Denote $\rho(M)$ by $\mu$, then $\rho(M^2)$ is $\mu^2$ and $R\infty\,(M^2) = 2\,R\infty\,(M)$. Thus the asymptotic rate of convergence of process (16) is twice that of process (4) for the linear case.

For the non-linear case, we consider (6) and (17). Let $\beta_i$ and $\gamma_i$ denote the eigenvalues of $I - G(X)$ and $I - 2G(X) + \overline{G}(X)\,G(X)$ respectively and as before let $\mu_i$ denote the eigenvalues of $I - G(Y)$. Let each eigenvalue be so indexed that $\beta_i$ is approximated by $\mu_i$ and $\gamma_i$ by $\mu_i^2$.

Now let $\sum a_i\mu^i, \sum b_i\beta^i$ and $\sum c_i\gamma^i$ represent the characteristic polynomials whose zeros are $\mu_i$, $\beta_i$ and $\gamma_i$ respectively. The coefficient $b_i$ is a continuous function of the elements of $G(X)$ and $|b_i - a_i|$ can be made as small as desired by choosing a proper sub-region $R_3$ in R. There is a subregion in $R_3$ having Y as its center such that $|\beta_i - \mu_i| < \epsilon$ provided X is in this region.

In a similar manner it can be shown that $|\gamma_i - \mu_i^2| < \epsilon$ by similarly choosing a proper subregion of $R_3$.

Now $\beta_i$ and $\gamma_i$ can be made arbitrarily close to $\mu_i$ and $\mu_i^2$ respectively by choosing the proper subregion of R. With the

condition that $\mu_i{}^2 < \mu_i$ it follows that the asymptotic rate of convergence for process (16) is double that of process (4).

## References

1. Bodewig, E. *Matrix Calculus*. North-Holland Publishing Co., Amsterdam, 1959.

2. Traub, J. F. *Iterative Methods for the Solution of Equations*. Prentice-Hall, Inc., Englewood Cliffs, 1964.

3. Hildebrand, F. B. *Introduction to Numerical Analysis*. McGraw-Hill Book Co., Inc., New York, 1956.

4. Varga, R. S. *Matrix Iterative Analysis*. Prentice-Hall, Inc., Englewood Cliffs, 1962.